

SciSpot

Paper # XXX

ABSTRACT

In this paper, we...

1 INTRODUCTION

Running parallel scientific applications, such as molecular dynamics (MD) simulations, on low-cost cloud transient resources.

The first "big" idea is that simulations are often bag of parallel tasks.

While there has been some past work that looks at running MPI applications on spot instances, our scope is much broader and considers how complete simulation pipelines can be run at low cost.

Spot on transient instances. Increasingly popular resource allocation model that is being offered by all cloud providers. Very low cost compared to conventional cloud resources, often by up to 10x. However, can be frequently revoked. Thus failure is a common occurrence, and not a rare-event. This is especially challenging for MPI jobs because of its inability to tolerate failures.

However, our insight is that while protecting a *single* job against revocations can require elaborate checkpointing based approaches, we don't necessarily have to do that if we consider that most simulations are composed of a series of jobs that search over a parameter space, and that what is important is the total running time and cost of this entire series of jobs.

Thus, no single job is "special".

Another aspect of novelty is that past work on transient resources used EC2 pricing information to get failure probabilities. However, this is no longer an accurate method. We perform the first empirical study of google preemptible VMs and their performance and availability for HPC workloads.

Another fundamental question is what is a suitable metric in such cases. Conventionally, it is speedup. In the cloud, it is some combination of cost and running time.

SciSpot is a framework and a tool that combines the use of failure modeling, checkpointing, and application-aware early stopping, to provide low cost execution of jobgroups for scientific applications.

Our work is the first to make a principled study of transient instances *other* than Amazon spot instances. Furthermore, our techniques make the first stab at addressing the new problems in the new EC2 spot pricing scheme.

Our work is in the context of reliability and cloud execution of scientific applications, and is novel because of multiple reasons: 1. Reliability and failure analysis of parallel scientific applications usually studied in the context of hardware with MTTFs of centuries, which is several orders of magnitudes higher than MTTFs faced in transient cloud servers (few hours).

2. While there have been studies of scientific applications been deployed in the context of cloud platforms, to the best of our knowledge, there has been no effort that integrates server selection and running jobgroups in a convenient automatic manner that makes it feasible to actually deploy applications on the cloud for scientists who may not have the requisite cloud experience.

2 BACKGROUND

2.1 Transient Computing

2.2 Parallel Scientific Applications in the Cloud

2.3 Molecular Dynamics Applications

Describe the kind of computation.

Scaling properties. Almost perfectly scalable with $O(n)$ communication?

This can be a like a case study of parallel scientific simulations. Will help relate to parameters etc with more concrete examples.

3 UNDERSTANDING AND MODELING TRANSIENT SERVER PREEMPTIONS

Transient cloud servers, by their very nature have limited availability and are frequently preempted. These preemptions are akin to fail-stop failures, and are often preceded by a small advance warning (few seconds) to allow for graceful shutdowns.

Since preemptions can impact the availability, performance, and cost of running applications, in this section, we examine their preemption characteristics. This modeling is important, because having a model of the availability can be useful in the context of predicting the running times of applications. Cloud providers offer a large number of servers of different configurations and types. Since transient server availability is fundamentally tied to supply and demand, the availability of servers of different types can be significantly different. Thus, selecting the "right" server type is crucial for minimizing the overall costs.

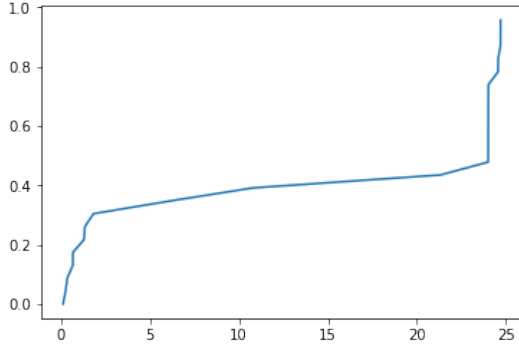


Figure 1: CDF of lifetimes of Google Preemptible Instances

3.1 EC2 spot instances

The earliest form of transient cloud instances. In addition to having dynamic availability, also have dynamic pricing. “Classic” spot instances had price determining the availability, and thus a large amount of work was devoted to bidding and analyzing the prices.

However a recent change to the spot prices no longer allows these assumptions, rendering it impossible to obtain the *exact* availability information from the prices alone.

3.2 Google Preemptible VMs

Launched in 2015. Flat-rate discount of 80% compared to on-demand servers. Interesting availability SLA: the maximum lifetime is 24 hours, and can be preempted earlier as well.

In this paper we will look at these preemptible VMs and show how to model their availability. Given the inability to use EC2 prices, we believe that our approach is more generalizable and robust.

There are some distinguishing characteristics of GCP preemptible VMs that makes their failure modeling challenging. First is their flat pricing and no other signalling information about their preemption rates (MTBFs) that makes server selection difficult.

3.3 Trade-offs in Server Selection

This presents us with many challenges in the cloud-deployment of these jobs.

Cloud providers offer multiple types of instances (VMs), with different hardware configuration (such as number of CPUs and memory size). The price of cloud servers is related to their hardware configuration, but it may not be strictly proportional to the hardware performance. For example, a VM with 32 CPUs may not be 32 times the cost of a single CPU VM.

For parallel and distributed applications, the type of servers selected has large implications on their performance. Consider the case of deploying an application on 8 8-core VMs vs. 16 4-core VMs. In both cases, the total number of CPU cores is the same. However, the larger number of VMs requires more communication between the application tasks, and thus may result in performance degradation. The performance of applications at different cluster configurations depends on their communication patterns and scaling properties.

Thus, when deploying applications on the cloud, one has to be mindful of the cost and performance tradeoff. However, in the case of transient servers, the story does not stop here.

In addition to pricing differences, the transient availability of instances *also* differs by type. Because the availability of a transient VM is broadly determined by the overall supply and demand of the instances of that *particular* type, the “preemption rate” of VMs often depends on the type of the instance.

Thus, selecting a transient cloud server involves a complex tradeoff between the cost of servers, their performance, and the preemption-rate. We develop server selection policies in the next section.

4 DESIGN

Scientific simulation applications consume a large amount of computational resources, and are often used in the context of *exploratory* research, where a large amount of jobs are run with different simulation parameters. This can be either a *parameter sweep*, where a large number of parameters need to be evaluated, or a *search* over a large parameter space for the “right” set of parameters that yield the desired model behavior.

In this paper, we look at the problem of running scientific simulations on *transient* computing resources in public clouds.

Past work has largely been focused on running parallel jobs (such as MPI) in the cloud. However, considering entire *job-groups* or ensembles of jobs presents new challenges and opportunities in timely, low-cost computation.

Our system, SciSpot, is a unified framework for running large job-groups that result from parameter exploration.

Input and some assumptions: We assume that the job-group consists of $J_1 \dots J_N$, with each job evaluating a model on some parameter. The list of parameters to explore can either be generated apriori (as in the case of parameter sweeps), or be dynamically generated as in the case of a search.

In this section, we will look at how we address these challenges:

- (1) How to select the right type of cloud server for an application?
- (2) How to effectively run job-groups?

SciSpot’s key insight is considering an entire job-group can allow better and simpler optimizations that can be easily deployed.

Job-groups are executed in two phases. In the first phase, we search for the right type of server for the jobs in the job-group, and then in the second phase, we execute the remaining jobs on the chosen servers.

4.1 Server Selection Policy

SciSpot’s server selection policy seeks to identify the best server type for a given job-group.

As stated in the previous subsection, the transient server selection problem is challenging because it involves balancing multiple optimization criteria: applications want low cost, low preemptions, and high performance.

Server selection based on application characteristics is a subject of a growing amount of recent work. These approaches often use micro benchmarks to gather performance data of cloud servers, and then use application performance models to determine suitable VMs for a specific application. Another class of approaches uses “black box” performance modeling, where the application’s performance is modeled using a function of the resources, for example, by using linear regression.

In contrast to prior work, our server selection employs a “cold start” policy, and we do not run profiling or pilot jobs that can increase the overall running time and cost. Instead, we search for the “best” cluster configuration for jobs in a job-group, by exploring the cluster configuration space for the “optimum” server type that optimizes all the desired parameters: cost, running time, and revocation rate.

Thus, the first e jobs in the job group $J_1 \dots J_e$ are the exploration jobs, run on different cluster configurations. We limit the total number of combinations to explore, by allowing users to submit an estimate of the total number of CPU cores that they desire for each job. This allows us to meet the user expectations in terms of performance and cost—whether the user expects us to spend a large amount of resources or not.

Thus, assuming that there are s different types of VM instances, the first s jobs are run on the s different types. Note that we use homogeneous clusters, since the performance of BSP programs in Heterogeneous environments can be degraded, and importantly, as we show, there are no performance or cost benefits to Heterogeneity.

For each server type i , we calculate the expected cost $E[C_i]$. $E[C_i] = n_i * c_i * E[T_i]$, where c_i is the price (per second) of the server, n_i is the number of servers of that type required to meet the core-count requirement. The expected running time of the job depends on two factors: the actual running time T , and the increase in running time due to preemptions. Each

preemption is akin to a fail-stop failure, which requires an application to restart. Our system makes no assumptions about the fault-tolerance policies supported by the application. For example, some applications may be able to *checkpoint* their state periodically. In either case (checkpointing or not), there is some work lost due to revocations. For ease of exposition, we assume no checkpointing. We discuss checkpointing in the next section (or never?)

Expected running time: Let the running time without failure be T :

$$E[T_i] = T + P(\text{at least one failure}) * T/2 \quad (1)$$

For calculating the probability of failure, we assume that the failure rate of an individual server of the type is p_i .

$$\begin{aligned} P(\text{at least one failure}) &= 1 - P(\text{no failure}) \\ &= 1 - (1 - p_i)^{n_i} \end{aligned} \quad (2)$$

Thus, we can see that if we select smaller VMs, we will require more of them (higher n_i), and this cluster configuration will have a larger probability of failure and thus higher running times and costs.

The probability of failure p_i depends on the type of server, and we use historically determined failure distributions. Roughly, if we assume exponentially distributed failures, then:

$$p_i = \frac{T}{\text{MTTF}_i} \quad (4)$$

Where MTTF is the mean time to failure of the server type, and T is the empirically determined job running time without failures (the best case).

In addition to searching over the servers, the effective number of servers is also dynamic in the case of transient environments due to preemptions. Thus, once we have found the appropriate server, we then explore the application’s performance at smaller cluster sizes, which helps in the job-group policies that we discuss next.

4.2 Preemption-handling Policies

We run the remaining jobs on the right configuration that is determined through the server selection policy.

Our goal is to minimize costs given a deadline.

This determines two things: how many jobs should be run in parallel, and what to do upon a revocation.

The number of jobs in parallel determines the overall size of our cluster.

$$\text{number of parallel jobs} = \frac{N \cdot T}{\text{Deadline Duration}}$$

If a server is preempted, then the job running on it will cease to run. Our preemption handling policies then decide what to do:

- (1) Restart the job on a smaller number of servers

- (2) Replenish lost servers and restart job
- (3) Discard job. This may be useful in case of parameter sweeps.

In addition, the user is also allowed to provide the fraction of jobs that are allowed to fail (η).

We exploit the naturally occurring intra-job redundancy.

New jobs. For new jobs, the questions are similar to the preemption policy ones, because they also must take into account the failures.

That is, assume that a job has finished and it is time to run a new job. So now, we have a set of servers that successfully ran a job. If we run on same set, then we may hit the 24 hour wall. But, if we discard these and launch new ones, then we face the infant mortality problem. So the question is, at what age should servers be retired? If they are too close to EOL (24 hours), then what's the point? Better start something fresh.

4.3 Checkpointing

Checkpointing requires the same number of servers, which may be tricky, whereas restarts can be on smaller number of nodes no problem.

4.4 Early Stopping

Based on the energy function, we can stop some simulations early. The early stopping criteria helps in minimizing the number of jobs run to completion.

We use this to proactively monitor jobs, as well as to decide whether to restart a job if it is preempted.

This can be a fairly substantial section

5 CHECKPOINTING POLICIES

High rate of failure means that especially for long jobs, checkpointing is necessary. In this section, we develop checkpointing policies.

At a high level, we only checkpoint long jobs, based on the probability of preemption.

Once checkpointing for a job is enabled, we use proactive periodic checkpointing.

DMTCP is used for checkpointing MPI programs.

Conventionally, periodic checkpointing according to Young-Daly formula: $\tau = \sqrt{2 \cdot \delta \cdot \text{MTTF}}$. This assumes that failure arrival process is exponentially distributed.

However, this may not always hold true. Clearly, from Figure 1, the CDF is not exponential.

Much more strongly "bath-tub". Maybe can use exponentiated Weibull to model the failures.

Regardless, because it is not memoryless, the checkpointing interval cannot be uniform. Maybe can use [2], which uses dynamic programming and also comes with a simulator of sorts.

6 IMPLEMENTATION

Central controller.

Cloud APIs for launching the jobs.

Slurm.

7 EVALUATION

The contenders:

- (1) Run every job on un-tuned on-demand instance (cost and running time)
- (2) Run every job on nanohub/big-red-2 (running and waiting time)
- (3) Run on transient, restart every time (cost)
- (4) SciSpot with early stopping and job sacrificing

7.1 Preemption likelihood curves

7.2 Searching for the best cloud configuration

7.3 Total cost vs. running time graphs

8 RELATED WORK

8.1 Scientific applications on cloud

A classic survey is [8] [16]

Parameter sweep: [3]

Price optimizations for Scientific workflows in the cloud [5]

8.2 Transiency mitigation

[9] classic work on MPI and Spot. Uses checkpointing. Redundancy, but for what? User specified number of VMs. Does not do instance selection. BCLR for checkpointing.

MOre spot and MPI: [6]. Focused on bidding and checkpoint interval. But bidding doesn't matter.

[12] is early work for spot and MPI and

[11] a batch computing service

Heterogeneity often used, but not useful in the context of MPI jobs. [10]

Selecting the best instance type, often for data analysis computations [1], and [15], and others like Ernest and Hemingway.

All the past work was on EC2 spot market with gang failures and independent markets [6, 9]. However this assumption has now changed, and failures can happen anytime. Our failure model is more general, and applies to both cases.

8.2.1 Fault-tolerance for MPI. [4] has a discussion of checkpointing frequency which is comprehensive.

Replication is another way [13]

8.2.2 Huge amount of work on bidding in HPC. [14]

[7]

8.3 Server Selection

Exploring a large configuration space using bayesian optimization methods in CherryPick [?] and Metis [?].

Can also use Latin Hypercube sampling for parameter exploration?

REFERENCES

- [1] ALIPOURFARD, O., AND YU, M. CherryPick: Adaptively Unearthing the Best Cloud Configurations for Big Data Analytics. 15.
- [2] BOUGERET, M., CASANOVA, H., RABIE, M., ROBERT, Y., AND VIVIEN, F. Checkpointing strategies for parallel jobs. In *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis on - SC '11* (Seattle, Washington, 2011), ACM Press, p. 1.
- [3] CASANOVA, H., LEGRAND, A., ZAGORODNOV, D., AND BERMAN, F. Heuristics for scheduling parameter sweep applications in grid environments. In *Proceedings 9th Heterogeneous Computing Workshop (HCW 2000) (Cat. No.PR00556)* (May 2000), pp. 349–363.
- [4] DONGARRA, J., HERAULT, T., AND ROBERT, Y. Fault tolerance techniques for high-performance computing. 66.
- [5] GARÁN, Y., MONGE, D. A., MATEOS, C., AND GARCÍA GARINO, C. Learning budget assignment policies for autoscaling scientific workflows in the cloud. *Cluster Computing* (Feb. 2019).
- [6] GONG, Y., HE, B., AND ZHOU, A. C. Monetary cost optimizations for MPI-based HPC applications on Amazon clouds: checkpoints and replicated execution. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis on - SC '15* (Austin, Texas, 2015), ACM Press, pp. 1–12.
- [7] GUO, W., CHEN, K., WU, Y., AND ZHENG, W. Bidding for Highly Available Services with Low Price in Spot Instance Market. In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing - HPDC '15* (Portland, Oregon, USA, 2015), ACM Press, pp. 191–202.
- [8] IOSUP, A., OSTERMANN, S., YIGITBASI, M. N., PRODAN, R., FAHRINGER, T., AND EPEMA, D. H. J. Performance Analysis of Cloud Computing Services for Many-Tasks Scientific Computing. *IEEE Transactions on Parallel and Distributed Systems* 22, 6 (June 2011), 931–945.
- [9] MARATHE, A., HARRIS, R., LOWENTHAL, D., DE SUPINSKI, B. R., ROUNTREE, B., AND SCHULZ, M. Exploiting redundancy for cost-effective, time-constrained execution of hpc applications on amazon ec2. In *HPDC* (2014), ACM.
- [10] SHARMA, P., IRWIN, D., AND SHENOY, P. Portfolio-driven resource management for transient cloud servers. In *Proceedings of ACM Measurement and Analysis of Computer Systems* (June 2017), vol. 1, p. 23.
- [11] SUBRAMANYA, S., GUO, T., SHARMA, P., IRWIN, D., AND SHENOY, P. SpotOn: A Batch Computing Service for the Spot Market. In *SOCC* (August 2015).
- [12] TAIFI, M., SHI, J. Y., AND KHREISHAH, A. SpotMPI: A Framework for Auction-Based HPC Computing Using Amazon Spot Instances. In *Algorithms and Architectures for Parallel Processing*, Y. Xiang, A. Cuzzocrea, M. Hobbs, and W. Zhou, Eds., vol. 7017. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011, pp. 109–120.
- [13] WALTERS, J. P., AND CHAUDHARY, V. Replication-Based Fault Tolerance for MPI Applications. *IEEE Transactions on Parallel and Distributed Systems* 20, 7 (July 2009), 997–1010.
- [14] WOLSKI, R., BREVIK, J., CHARD, R., AND CHARD, K. Probabilistic guarantees of execution duration for Amazon spot instances. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis on - SC '17* (Denver, Colorado, 2017), ACM Press, pp. 1–11.
- [15] YADWADKAR, N. J., HARIHARAN, B., GONZALEZ, J. E., SMITH, B., AND KATZ, R. H. Selecting the *best* VM across multiple public clouds: a data-driven performance modeling approach. In *Proceedings of the 2017 Symposium on Cloud Computing - SoCC '17* (Santa Clara, California, 2017), ACM Press, pp. 452–465.
- [16] ZHAI, Y., LIU, M., ZHAI, J., MA, X., AND CHEN, W. Cloud versus in-house cluster: evaluating Amazon cluster compute instances for running MPI applications. In *State of the Practice Reports on - SC '11* (Seattle, Washington, 2011), ACM Press, p. 1.